

Predicting English Premier League Games from Match Statistics

Machine Learning Final Project Part 2

Bhavik Naik

Faculty of Business and I.T

Ontario Tech University

Oshawa Ontario Canada

bhavik.naik@ontariotechu.net

1 Frame the Problem

Sports is one of the most popular hobbies in the world. Athletes train daily to stay ahead of their competition while fans enjoy watching the action unfold on live TV. The Premier League in England is the most watched football league in the world [SportMob, 2021]. Players from around the world compete against 20 teams to be claimed the best English football side in the country. With 925.3 million live streaming minutes and an average of 462,000 viewers per match in the 2019/20 season at NBC Sports Network, the Premier League's fan base extends further than England itself [SportsPro, 2020]. Super computers during the beginning and middle of the season always makes headlines as these models predict the final league standings based on the information that has already been collected. Nowadays, there is a lot of data collected per game. From the basic numbers such as shots on goal and kilometers run, to more advanced numbers such as the expected goal (xG) which takes in player positioning and other data to calculate the chance of goal that the team has. Teams and managers have altered their playing styles to improve this number so that the possibility of a goal scoring chance is higher. By recording all these data points, we can utilize regression algorithms to calculate the results based on the key factors that influence the game. If we are successfully able to find a model that can take in data points and accurately calculate the match result, we can theoretically predict what could happen in future games. The dataset that was chosen for this paper is the "English Premier League stats 2019-2020" which contain 45 columns of data corresponding to each game, respective to a single team. The xG, scored, conceded, pressing plays, number of shots, betting odds, and more are contained in this dataset. In this paper, we will train three machine learning models to take in this data and predict the correct results. We will be using a deep neural network, linear regression, and SGD regressor as our three machine learning models. We will explore the data first, clean up the data, and finally, showcase the results from all three. The best performing algorithm will get graphs depicting the results.

2 Gain Insights

In this section, we will provide four graphs that try to gain insights into the data.

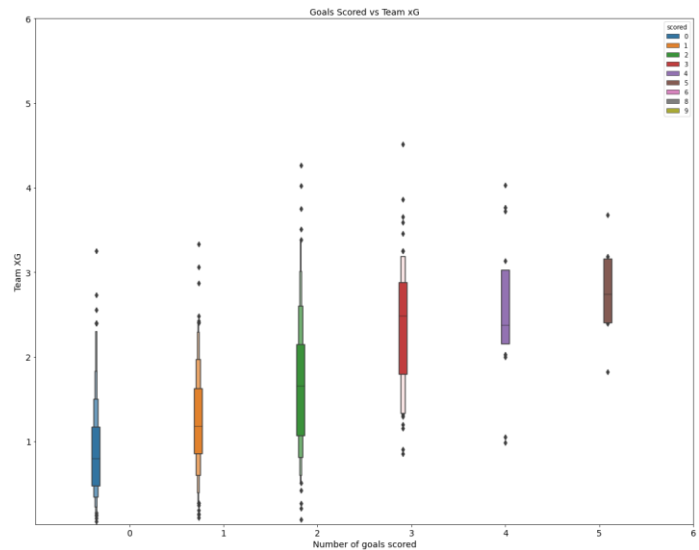


Figure 1 The correlation between goals scored (x-axis) and the team xG. We can see that the higher the xG, the more goals are being scored by the team.

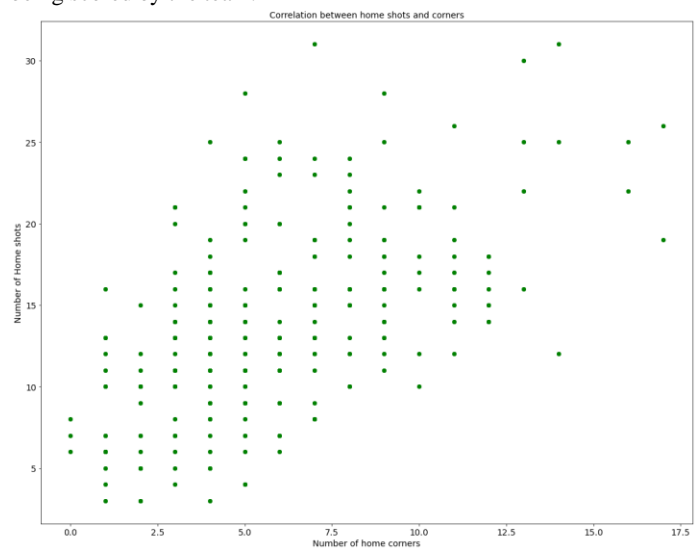


Figure 2 The number of home corners (x-axis) vs the number of shots. As we can see, the more corners a team has, the more shots they have on average. It is quite interesting to see correlations in areas that you would not expect.

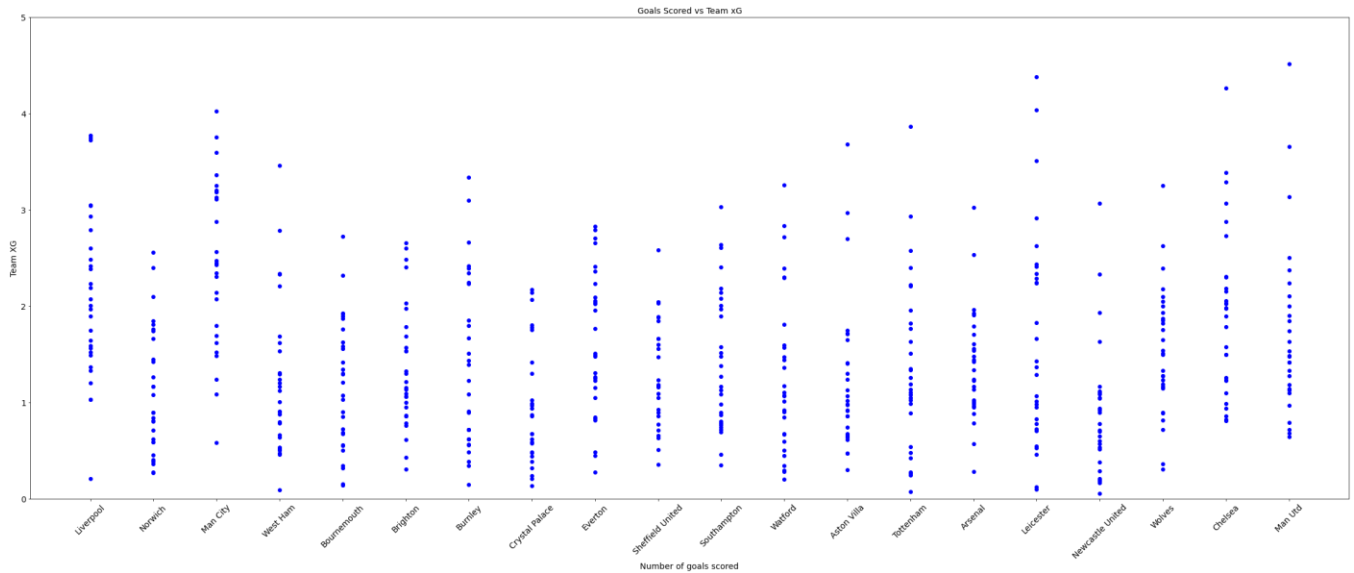


Figure 3 A scatter plot showing the team's xG scores across their matches. We can see that higher ranked teams like Liverpool, Chelsea, and Manchester City have higher xG scores because they attack more than lower ranked teams like Newcastle United.

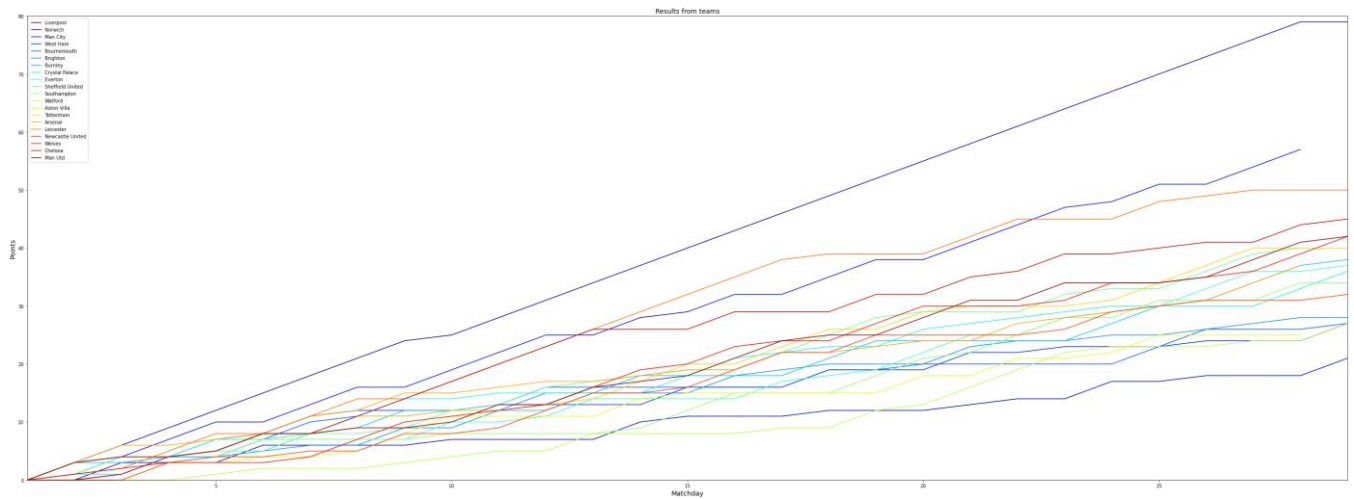


Figure 4 The results from each game, mapped over the entire season. We can see every team and how they were able to accumulate points throughout the season.

3 Prepare the Data

The information that was present in the dataset did not include any null values which prevented us to remove lines that contained missing data. There were numerous columns that we had to drop which did not make sense in calculating the result, such as the referee name, matchday (the day that the match was on), date, match time, round (the number of played matches by the team in the season), points (the number of points that team gained after the match), total points (the number of points that team has coming into the match). All the bet odds were also dropped which we felt did not contribute to the result of the game. There were some columns that were dropped that included numbers which could help the algorithm, but we felt that these numbers did not further enhance the result of the algorithm. Those include, npxG and npxA (the expected goal value without penalties, where each penalty adds 0.76 to the XG), and npxD (the difference XG between the teams). In terms of the data, we had to convert some of the columns to integers, such as the home and away column,

and the result column. The values were replaced from *h* (home) and *a* (away) to 0 and 1. While in the results column, *w* (win) was represented as 2, *d* (draw) was represented as 1, and *l* (lost) was represented as 0. We also converted the 'teamId' field which contains the team's name to integers. Each team got a number from 0 to 19 and this was inserted into the data frame in expense of the name of the team. We decided to keep in team name because higher ranked teams perform better than their counterparts. We then used the SKLearn `train_test_split` function to split the data into training and testing tests using a 60-40% split. We set the test column to the results column as this is the column that we are trying to predict.

4 Machine Learning Algorithms

In this paper, we used three machine learning algorithms to try to predict the results based on the data in the column. We will be looking at a deep neural network (DNN) that is based of Keras, a linear regression algorithm, and a Stochastic gradient descent (SGD) regressor to evaluate the data and predict results. We will

be comparing the three algorithms using mean squared error and mean absolute error. Table 1 will showcase the comparative performance table which will contain the results from each algorithm.

4.1 Deep Neural Network

Starting off with the deep learning algorithm, we created a Keras sequential model that uses seven layers. It starts off with a flatten layer that takes in the 29 columns, then a 128-node connected layer with activation set to “relu” and kernel initializer set to normal. The third layer is a dropout layer with 30% of nodes dropping out. Then came another 128-node connected layer with the kernel initializer set to normal and activation set to “relu”. A drop out layer with a 50% chance came after and the sixth layer had a 64-node connected dense layer with kernel initializer set to normal and activation set to “relu”. Finally, the last layer was a 3-node dense with the activation set to linear. We used mean squared error as the loss and the metrics we added include the mean squared error and the mean absolute error. The optimizer we used was “adam”. When training it, we set epochs to 100 and fit it with the training set while the validation set was set to the testing set.

4.2 Linear Regression

The second algorithm that we used to predict the results was a linear regression algorithm. A linear regression algorithm is one of the most popular machine learning models to use when dealing with regression problems. We fit the testing data on it and left the parameters as the default depending on what is given by Scikit-learn.

4.3 SGD Regressor

The final machine learning model that we used is a Stochastic gradient descent or SGD regressor. The SGD regressor is also one of the popular regression algorithms and it is tailored to a small amount of data. We used the SkLearn StandardScaler function to scale up the testing data, and then, we called the SGDRegressor with all its default values. After training it on the training dataset, we predicted the results of the testing dataset.

4.4 Results

We will now discuss the results from the three algorithms.

4.4.1 Deep Neural Network. For the results, the DNN gave us a score of 0.0150 for the mean squared error while the mean absolute error was 0.0852. The lower the better when it comes to these two metrics, and we can see that the deep neural network performed fairly well as it was able to keep both metrics under 1.

4.4.2 Linear Regression. The linear regression performed much better than the deep learning model. The mean squared error was 1.09×10^{-29} and the mean absolute error was 2.58×10^{-15} . This shows that the algorithm got close to 0, which means that it was able to

predict the result nearly precisely from the test data that it received.

4.4.3 SGD Regressor. The SGD regressor gave us a mean squared error of 0.0024037 and a mean absolute error of 0.0340784. It also performed well as both values are significantly below 1.

Algorithm	Mean Squared Error (MSE)	Mean Absolute Error (MAE)
Deep Neural Network	0.0150	0.0852
Linear Regression	1.09×10^{-29}	2.58×10^{-15}
SDG Regressor	0.0024037	0.0340784

Table 1: The comparative performance table

As we can see from table 1, the linear regression algorithm performed the best and got the closest to predicting the testing data. The SDG regressor was second and then the deep neural network performed the worst.

5 Best Performing Algorithm

In this section, we will show four graphs that showcase the results of the linear regression model as this is the best performing algorithm.

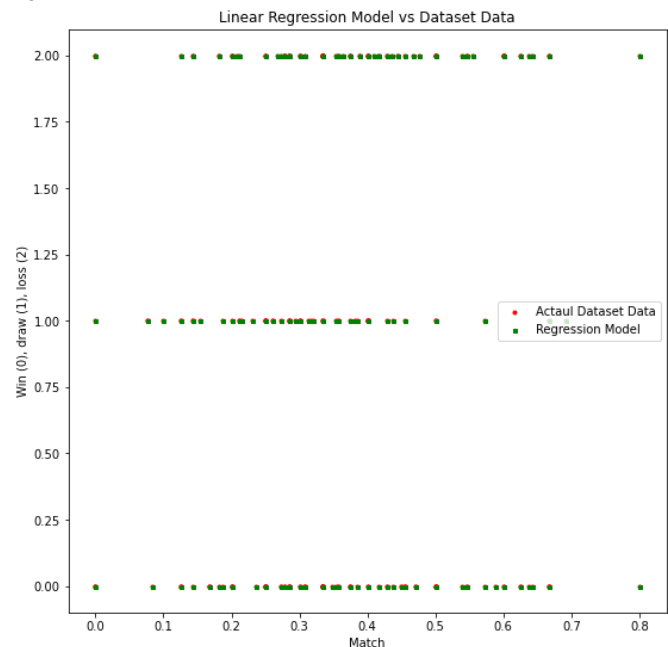


Figure 5 The first graph shows a scatter plot of the testing data (red) vs the regression model (green). As we can see, the points are similar, almost overlapping. It is difficult to see the red dots as the green ones are directly on top of them. On the y axis, 0 represents a loss, 1 represents a draw and 2 represents a win

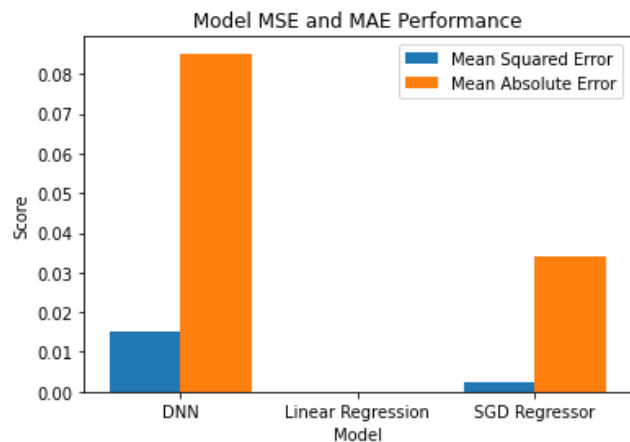


Figure 6 The MSE and MAE performance across the models. As we can see, the linear regression model did so well that it does not fit in the scale when comparing the other 2 models

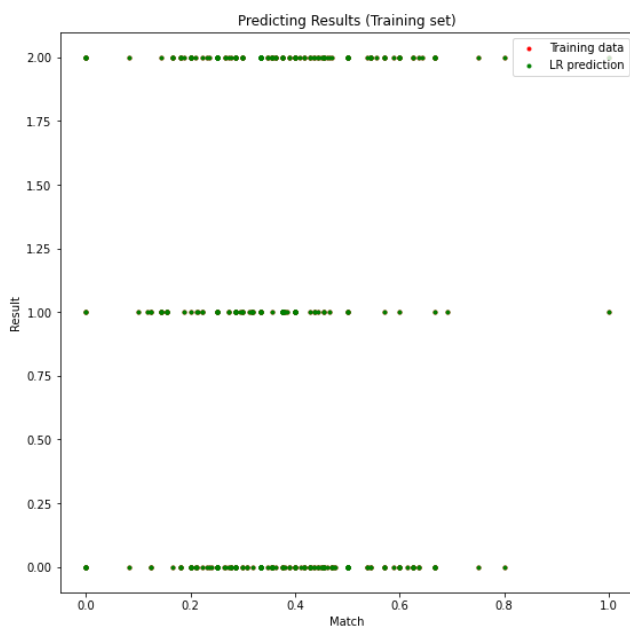


Figure 7 This figure is similar to figure 5 but instead of showing the performance of the linear regression model in terms of the testing data, we look at how well it can predict the training data. Based on this graph, we can see that it easily is able to predict the training data too.

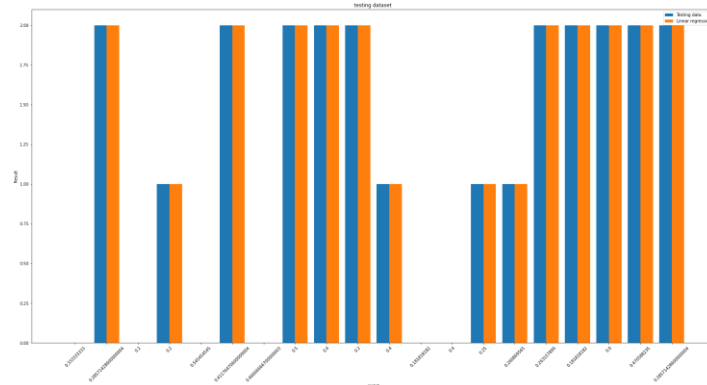


Figure 8 This bar graph shows how similar the linear regression model is to the testing data. We can see that the orange bar (linear regression) is exactly the same as the testing data. This is just a slice of the data as graphing all of the data will not fit.

APPENDIX

In the appendix, we will show all of the Python code that we used to achieve the results of the report.

```
1. # Bhavik Naik
2. # Machine learning final project
3.
4. # Import the various libraries
5. import pandas as pd
6. import os
7. import math
8. import numpy as np
9. import tensorflow as tf
10. from tensorflow import keras
11. import keras.optimizers
12. import numpy as np
13. import seaborn as sns
14. import matplotlib.pyplot as plt
15. import matplotlib.cm as cm
16. import plotly.graph_objects as go
17. import plotly.express as px
18. from google.colab import files
19. import io
20. uploaded = files.upload()
21. # link to the kaggle
22. #https://www.kaggle.com/idooyo92/epl-stats-20192020
23. # read in the csv
24. df = pd.read_csv(io.BytesIO(uploaded['epl2020.csv']), index_col=0)
25.
26. # replace the home and away to numbers, and win, lost and draw
27. df = df.replace({'h': 0, 'a': 1})
28. df = df.replace({'w': 2, 'd': 1, 'l': 0})
29. df.head()
30. # create a dictionary with the team results over the season to graph it (figure 4)
31. teamresults = {}
32. teamresults2 = {}
33. for i in df['teamId']:
34.     teamresults[i] = list((df.loc[df['teamId'] == i, 'pts']))
35. for i in teamresults:
36.     teamresults2[i] = []
37. for i in teamresults:
38.     for j in range(len(teamresults[i])):
39.         teamresults2[i].append(sum(teamresults[i][:j]))
40. # plot the lines on the graph, with the results
41. f1 = plt.figure(figsize=(50, 18))
42. # set the colour
43. plt.gca().set_prop_cycle(plt.cycler('color', plt.cm.jet(np.linspace(0, 1, 20))))
44. labels = []
45. for i in teamresults2:
46.     plt.plot(range(1, len(teamresults2[i]) + 1), teamresults2[i])
47.     labels.append(i)
48. plt.legend(labels, loc="upper left")
49. plt.xlim([1, 29])
50. plt.ylim([0, 80])
51. plt.xlabel("Matchday", fontsize="14")
52. plt.ylabel("Points", fontsize="14")
53. plt.title("Results from teams", fontsize="14")
54.
55. # create the graph for the teamid and the xg for the team (figure 3)
56. f3 = plt.figure(figsize=(40, 15))
```

```

57. plt.scatter(x = df['teamId'] , y = df['xG'], color="blue")
58. plt.xticks(fontsize=14, rotation=45)
59. plt.yticks(fontsize=14)
60. plt.xlabel(xlabel="Number of goals scored", fontsize=14)
61. plt.ylabel(ylabel="Team xG", fontsize=14)
62. plt.title("Goals Scored vs Team xG", fontsize=14)
63. plt.ylim([0, 5])
64.
65. # replace the teamid names to integers. use a dictionary to do this
66. a = 0
67. values = {}
68. for i in df['teamId']:
69.     if i not in values:
70.         values[i] = a
71.         a += 1
72. # replace the teamid with the integers
73. df = df.replace(values)
74.
75. # drop all of the axis we dont need
76. df = df.drop(["Referee.x", "matchDay", "date", "matchtime", "pts", "tot_points", "round",
77.               "matchtime", "npGD", "xpts", "npG", "npGA",
78.               "B36SH.x", "B36SD.x", "B36SA.x"], axis=1)
79.
80. import seaborn as sns
81.
82. # create the correlation graph between home corners and home shots (figure 2)
83. f2 = plt.figure(figsize=(19, 15))
84. plt.xticks(fontsize=14)
85. plt.yticks(fontsize=14)
86. plt.scatter(x=df["HC.x"], y=df["HS.x"], color="green")
87. plt.title("Correlation between home shots and corners", fontsize=14)
88. plt.xlabel("Number of home corners", fontsize="14")
89. plt.ylabel("Number of Home shots", fontsize="14")
90. plt.show()
91.
92. # create the graph showing the correlation between goals scored and the team xg (figure 1)
93. f3 = plt.figure(figsize=(19, 15))
94. ax = sns.boxenplot(data = df, x = df['scored'] , y = df['xG'], hue=df['scored'])
95. plt.xticks(fontsize=14)
96. plt.yticks(fontsize=14)
97. plt.xlabel(xlabel="Number of goals scored", fontsize=14)
98. plt.ylabel(ylabel="Team xG", fontsize=14)
99. plt.title("Goals Scored vs Team xG", fontsize=14)
100. plt.xlim([-1, 6])
101. plt.ylim([0.02, 6])
102.
103.
104. # create the training and testing data
105. from sklearn.model_selection import train_test_split
106. train, test = train_test_split(df, test_size=0.4)
107. train_labels = train.iloc[:, -1]
108. train_data = train.drop(['result'], axis=1)
109. test_labels = test.iloc[:, -1]
110. test_data = test.drop(['result'], axis=1)
111. train_labels_mc = train['result']
112. test_labels_mc = test['result']
113.
114. # Create the Deep Neural Network and build the model with the nodes
115. model = keras.models.Sequential([
116.     keras.layers.Flatten(input_shape=[29]),

```

```
117. keras.layers.Dense(128, activation="relu", kernel_initializer='normal'),
118. keras.layers.Dropout(rate=0.3),
119. keras.layers.Dense(128, kernel_initializer='normal',activation='relu'),
120. keras.layers.Dropout(rate=0.5),
121. keras.layers.Dense(64, kernel_initializer='normal',activation='relu'),
122. keras.layers.Dense(3, activation='linear')
123. ])
124. model.summary()
125.
126. # Compile the model with the correct loss, optimizers, and metrics
127. import keras.optimizers
128. model.compile(loss="mean_squared_error",
129. optimizer="adam",
130. metrics=["mean_squared_error", "mean_absolute_error"])
131.
132. # fit the training and testing datasets to the dnn model
133. history = model.fit(train_data, train_labels_mc, epochs=100, validation_data=(test_data, test_labels_mc), verbose=1)
134.
135. from sklearn.linear_model import LinearRegression
136. from sklearn.preprocessing import StandardScaler
137.
138. # create the linear regression model and fit the data
139. reg = LinearRegression()
140. reg.fit(train_data, train_labels_mc)
141. # predict the testing set
142. matchResultLR = reg.predict(test_data)
143. from sklearn.metrics import mean_squared_error, mean_absolute_error
144. # print out the mean squared error and mae
145. print(mean_squared_error(test_labels_mc, matchResultLR))
146. print(mean_absolute_error(test_labels_mc, matchResultLR))
147.
148. from sklearn.pipeline import make_pipeline
149. from sklearn.preprocessing import StandardScaler
150. from sklearn.linear_model import SGDRegressor
151. # create the sgd regressor model
152. sgdreg = make_pipeline(StandardScaler(), SGDRegressor())
153. # fit the data
154. sgdreg.fit(train_data, train_labels_mc)
155. # predict the data
156. sgd = sgdreg.predict(test_data)#
157. # print out the mse and mae
158. print(mean_squared_error(test_labels_mc, sgd))
159. print(mean_absolute_error(test_labels_mc, sgd))
160.
161. # create the graph to show the performance of the regression model (figure 5)
162. fig = plt.figure(figsize = (9, 9))
163. ax1 = fig.add_subplot(111)
164. # Create the actual scatter plot
165. ax1.scatter(test_labels, test_labels_mc, s=9, c='red', marker="o",
166. label='Actual Dataset Data')
167. ax1.scatter(test_labels, matchResultLR, s=9, c='green', marker="s",
168. label='Regression Model')
169. plt.legend(loc='upper right')
170. plt.xlabel('Match')
171. plt.ylabel('Win (0), draw (1), loss (2)')
172. plt.title('Linear Regression Model vs Dataset Data')
173. plt.legend()
174. plt.show()
175.
176. # create the graph to show the bar graph with the 3 models comparison (figure 6)
```

```
177. X = ['DNN', 'Linear Regression', 'SGD Regressor']
178. Y = [0.0150, 1.091e-29, 0.0024037]
179. Z = [0.0852, 2.58e-15, 0.0340784]
180. X_axis = np.arange(len(X))
181. # make the bar graph
182. plt.bar(X_axis - 0.2, Y, 0.4, label = 'Mean Squared Error')
183. plt.bar(X_axis + 0.2, Z, 0.4, label = 'Mean Absolute Error')
184. plt.xticks(X_axis, X)
185. plt.xlabel("Model")
186. plt.ylabel("Score")
187. plt.title("Model MSE and MAE Performance")
188. plt.legend()
189. plt.show()
190.
191. # make the chart showing the performance of the linear regression model with the training set
192. fig = plt.figure(figsize = (9, 9))
193. plt.scatter(train_labels, train_labels_mc, color = "red", label="Training data", s=9)
194. plt.scatter(train_labels, reg.predict(train_data), color = "green", label="LR prediction", s=9)
195. plt.title("Predicting Results (Training set)")
196. plt.xlabel("Match")
197. plt.ylabel("Result")
198. plt.legend(loc='upper right')
199. plt.show()
200.
201. # make the bar graph showing the performace of the linear regression model (figure 8)
202. X = test_labels[10:30]
203. Y = test_labels_mc[10:30]
204. Z = matchResultLR[10:30]
205. X_axis = np.arange(len(X))
206. fig8 = plt.figure(figsize = (30, 15))
207. # make the bar graph
208. plt.bar(X_axis - 0.2, Y, 0.4, label = 'Testing data')
209. plt.bar(X_axis + 0.2, Z, 0.4, label = 'Linear regression')
210. plt.xticks(X_axis, X, rotation=45)
211. plt.xlabel("match")
212. plt.ylabel("Result")
213. plt.title("testing dataset ")
214. plt.legend()
215. plt.show()
```